



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

Bandwidth-Sharing Schemes for Multiple Multi-Party Sessions

Zhen Liu — Naceur Malouch — Vishal Misra — Dan Rubenstein [‡] — Sambit Sahu ^{*}

N° 4821

May 2003

THÈME 1

A large blue rectangle occupies the lower half of the page. Overlaid on it is a large, light gray stylized 'R' that is partially cut off on the left. To the right of the 'R', the words 'Rapport de recherche' are written in a white serif font. A horizontal gray brushstroke is positioned below the text.

*Rapport
de recherche*



Bandwidth-Sharing Schemes for Multiple Multi-Party Sessions

Zhen Liu ^{*}, Naceur Malouch [†], Vishal Misra [‡], Dan Rubenstein [‡], Sambit Sahu ^{*}

Thème 1 — Réseaux et systèmes
Projet Mistral

Rapport de recherche n° 4821 — May 2003 — 26 pages

Abstract: A recent means for enabling multicast in the Internet involves deploying network overlays where end-systems participate in the forwarding of data to other end-systems. The use of overlays not only permits individual sessions to simultaneously structure their communication trees atop unicast-only networks, but also gives the session greater flexibility when forming the topology of the forwarding tree. Even though multiple sessions are often expected to compete for the same network overlay resources, most work to date assumes that overlay protocols operate as though each session has isolated access to the available overlay resources. We consider two algorithms that build depth-bounded overlay trees where each node's outgoing bandwidth constrains the number of nodes to which it can directly forward data. One algorithm tries to cluster a node's available bandwidth within a single tree, the other tries to disperse the available bandwidth among multiple trees. We prove analytically that when node capacities are identical and session requirements are identical that a clustering approach will increase the number of sessions that can co-exist. However, simulation results reveal that in heterogeneous networking environments or in environments where session participants vary with time, the dispersing algorithm outperforms the clustering algorithm. These results can be used to guide future development of overlay protocols that must partition resources among multiple sessions.

Key-words: Multicast, Application Layer, Network Overlays, Optimization Problem, Delay Constraint, Majorization, Heuristics, Routing

^{*} IBM T.J. Watson Research Center, P.O.Box 704, Yorktown Heights, NY 10598, USA

[†] INRIA Sophia-Antipolis, 2004 Route des Lucioles, Sophia-Antipolis 06902, France

[‡] Columbia University, New York, NY 10027, USA

Partage de bande passante en multicast applicatif

Résumé : Le principe du multicast applicatif est de construire un arbre multicast, reliant les participants d'une session multicast, en utilisant les connexions au niveau transport déjà offertes par les protocoles Internet actuels. En d'autres termes, le réseau ne joue aucun rôle particulier pour transmettre les données multicast, les données sont transmises moyennant le routage unicast normal. Contrairement à l'IP-Multicast, plusieurs copies des données peuvent circuler sur le même lien physique. En outre, le nombre de connexions générées par chaque noeud (participant) est limité par sa bande passante. Cette limitation est encore plus stricte pour les machines clientes qui utilisent en général les modems ou le DSL pour se connecter au réseau. La contrainte de bande passante peut aussi engendrer la limitation du nombre de participants et un délai très élevé entre la source et certains participants. L'optimisation de l'usage de bande passante est une question cruciale dans le multicast applicatif parce que souvent plusieurs sessions multicast concourent pour les mêmes ressources du réseau. Cependant, la plupart des protocoles supposent implicitement qu'elles ont un accès isolé aux bandes passantes. Nous avons proposé deux nouveaux algorithmes qui visent à maximiser le nombre de sessions de multicast applicatif qui peuvent coexister dans le même réseau. Nous prouvons analytiquement que quand les capacités des nœuds et aussi les caractéristiques des sessions sont identiques, un algorithme "groupant" est optimal. Dans un environnement plus hétérogène, nous employons des simulations pour étudier l'impact des algorithmes sur la probabilité de blocage dans divers scénarios comprenant des sessions à population variable où l'adhésion des membres est dynamique. Nous modélisons également l'environnement hétérogène par un système de files d'attente grâce auquel nous avons déterminé une limite inférieure théorique sur la probabilité de blocage qui s'applique à toutes les solutions de partage de bande passante. Nous avons trouvé que la politique "dispersante" est la plus adaptée dans la plupart des cas. Ces résultats peuvent être employés pour aider le développement des protocoles de recouvrement qui cherchent à répartir les ressources parmi des sessions multiples.

Mots-clés : Multicast applicatif, réseaux de recouvrement, problème d'optimisation, contrainte de délai, majoration, heuristique, routage

1 Introduction

Multi-party applications such as distributed gaming, tele-video conferencing, and distributed concerts require simultaneous transmission from a set of source points to possibly overlapping sets of receiving points. For instance, in a distributed game, to keep all players' perspectives consistent, each player must update their position information and status to other players whose virtual positions are located close by [8]. Tele-video conferencing and distributed concerts require that the participants can almost simultaneously see and hear one another. Hence, it is important to configure these sessions such that transmissions are received within a given delay bound.

A viable approach today to provide multipoint communication over the IP network is via a *multicast overlay* [7, 6, 5]: a set of end-systems that connect the source of the transmission to all receivers in the form of a tree. Each intermediate node on the tree takes transmissions sent from its parent node and forwards these transmissions to its children nodes within the tree. The forwarding itself is often handled via network and transport level unicast protocols such as TCP, and UDP. While there have been several recent works [7, 6, 5] that develop efficient algorithms for building overlays for group communication, few have explored how to construct session topologies in environments where multiple sessions, whose memberships vary dynamically with time, compete for the same network resources. Previous work in multicast that considered multiple sessions competing for resources is in the context of congestion control [4, 13, 20, 15]. There, sessions are expected to reduce bandwidth requirements when insufficient bandwidth is available to avoid proliferating a congested state within the network. Bounded-fanout multicast has been explored in the context of building (non-source-specific) minimum spanning trees that minimize aggregate edge costs instead of minimizing the delay from a specific source [3, 2]. [19] presents a survey of previous work in the area of QoS multicast routing. However, there is no discussion in the survey of work that addresses the problem we consider here.

In this paper, we evaluate overlay construction algorithms that construct overlay trees in environments where each node bounds its fanout: the number of neighbors to which it can forward transmissions. Fanout bounds are often necessary in such overlay networks, since the end-systems that act as forwarding agents utilize modem, DSL or cable links and can only offer limited bandwidth to the sessions they support. Our optimization criteria is to build *depth-bounded* trees such that the number of hops through the overlay that are taken to reach a receiver within a certain bound. An algorithm that minimizes a tree's depth does not necessarily minimize the end-

to-end propagation delay to a receiver as the delays between different pairs of overlay nodes can differ. Nonetheless, there are several reasons why depth-bounded trees are of interest. First, it is a non-trivial task to obtain accurate estimates of propagation delays between overlay participants. Second, delays that occur in transmission between such end-systems is often dominated by a combination of the delay across this last-mile technology and the lack of priority given to processing of arriving and departing packet transmissions.

In this paper, we consider two baseline algorithms, CLUSTER and DISPERSE, whose approaches toward selecting the nodes that support transmission for a session lie at two extremes. CLUSTER attempts to apply a node's outgoing bandwidth toward a single session, thereby *minimizing* the number of sessions for which a node acts as a forwarding agent. In contrast, DISPERSE splits a node's outgoing bandwidth across as many sessions as possible while maintaining the depth-bound in the constructed tree, thereby *maximizing* the number of sessions for which a node acts as a forwarding agent.

We begin by showing via a mathematical analysis that CLUSTER is optimal in homogeneous environments where all nodes have the same bandwidth limitations and all multicast sessions have the same set of requirements. We then turn our attention to heterogeneous networking environments where nodes have variable capacity, and overlay participants can join and leave the network. We begin by modeling the heterogeneous environment as a queueing system and use this system to derive a theoretical lower bound on the blocking probability of algorithms in these settings. This lower bound translates into an upper bound on the rate at which sessions can be admitted into the network. Since no algorithm can perform better than this theoretical bound, it provides an additional benchmark to evaluate algorithm performance. Here, we show near optimality by demonstrating that, by merely allowing the sender of an arriving session to transfer its current transmission responsibilities for other sessions to other nodes, CLUSTER and DISPERSE approach the theoretical bound. Last, we use simulation to evaluate these algorithms in settings where session membership varies with time. We find that, although CLUSTER is optimal in homogeneous settings where session participation is static, DISPERSE either performs as well or outperforms CLUSTER, yielding lower blocking rates than CLUSTER in heterogeneous settings or where session membership varies.

Our work builds on the recent work such as [7, 6, 5, 14, 10, 1, 9] whose focus was producing “good” overlays for a single source's multicast transmissions. None of these works explicitly examines simultaneously supporting multiple multicast sessions or sources using a shared pool of overlay resources. This is not due to an expectation

that only one source's transmissions will utilize overlay resources at any given time, but because the heterogeneity, dynamics, and diversity in multicast session requirements and network capabilities make finding “good” practical solutions at this point quite difficult. Recent investigation into a similar problem has been considered in [17, 18]. Our work not only presents a contrasting solution, but our solution and evaluation goes further by also capturing the important dynamics of participants joining and leaving a session when the session is already and continues to be active. Furthermore, our theoretical results provide performance bounds that apply to any solution in the area.

The rest of the paper proceeds as follows. In Section 2, we introduce the general model of the network and multicast sessions for our analysis. Section 3 describes the various algorithms that we consider. In Section 4, we demonstrate optimality of one of the algorithms in a homogeneous networking environment. In Section 5, we introduce a Stochastic Knapsack queuing model that we use to compute lower bounds on blocking probabilities in heterogeneous networking environments. In Section 6, we compare the performance of our various proposed algorithms to one another and to the lower bound. In Section 7 we compare the performance of the algorithms when session members dynamically leave and join, and Section 8 concludes the paper.

2 Network Model

In this section, we present the network and session models that will be used to evaluate our algorithms that form session trees for multiple sessions. Let \mathcal{N} be the set of receivers, i.e., the set of nodes that wish to participate in at least one multicast overlay session. We define $\mathcal{S} \subset \mathcal{N}$ to be a sequence of these nodes that wish to act as transmitters for a session. Note that this set can contain repetitions, i.e., a node appears k times in \mathcal{S} if it is the sender for k different sessions. We denote the i th entry in the sequence by s_i .

Let $\mathcal{R}_i \subset \mathcal{N}$ be a set of nodes that wish to receive the transmissions from sender s_i . By convention, we require that $s_i \in \mathcal{R}_i$. We assume that a node n is only willing to participate in overlay multicast session when it is interested in receiving the transmitted data of that session. In other words, \mathcal{R}_i is exactly the set of nodes that form the multicast tree that will carry data for the session for which s is the sender. Note that there will be a set of leaf nodes in this multicast tree that will be receiving, but will not be forwarding these transmissions.

The i th session forwards data to all participants at an arbitrary rate, ρ_{s_i} . β_n is the aggregate outgoing bandwidth of node n , such that if n forwards directly to c_j nodes

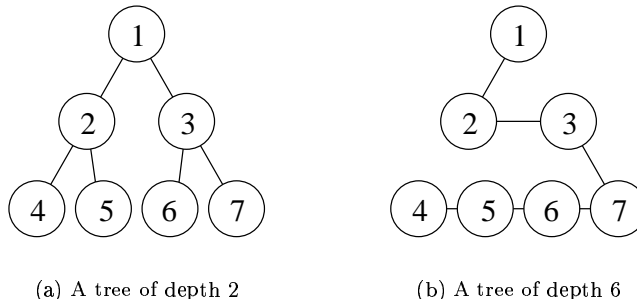


Figure 1: Selecting “good” clusterings

for session j , then it must be the case that $\beta_n \geq \sum_{j \in A} c_j \rho_{s_j}$ where A is the set of active sessions. This requirement is applicable to several last-mile technologies such as cable and DSL that provide asymmetric bandwidths in the two directions, where incoming bandwidth is significantly larger (and essentially unbounded) in comparison to the outgoing bandwidth.

We define Δ_i to be a bound on the depth of the tree for the session whose transmissions emanate from s_i , i.e., no node should be more than Δ_i overlay hops from the source node. For simplicity of presentation, we assume that $\Delta = \Delta_i$ is constant for all sessions.

3 Algorithms

We now describe the algorithms that we use to determine whether or not to admit an incoming session, and, when admitted, the topology structure of the tree.

To give the reader a feel for what makes a “good” algorithm, consider the trees constructed in Figure 1. Assume that each node in this figure is capable of providing direct transmissions to two additional nodes. In Figure 1(a), 3 nodes have been selected, each of which is fully utilizing its outgoing transmissions to connect all nodes within a tree in which node 1 is the source. The tree formed has depth 2, and, in addition, another tree of depth 2 can be constructed (e.g., node 4 transmits to nodes 5 and 6, node 5 transmits to nodes 2 and 3, node 6 transmits to nodes 7 and 1). Hence, trees that cluster nodes within a session are intuitively “good”. In contrast, the tree in Figure 1(b) has depth 6 and precludes any other trees of depth

less than 6 from being constructed. This figure demonstrates that intuitively, it is preferable to utilize the outgoing edges of a node within a single tree, rather than to spread these edges among several trees. We shall see, however, that this intuition can be misleading in dynamic and heterogeneous environments.

3.1 Algorithm CLUSTER

CLUSTER forms trees by utilizing a set of nodes for forwarding such that the nodes' remaining available bandwidth is minimized. The algorithm distinguishes between *partial nodes*: nodes whose bandwidth has already been applied within other sessions, and *full nodes*: nodes that do not currently forward transmissions on behalf of other sessions. CLUSTER tries to incorporate as many already-existing partial nodes as possible into the middle (i.e., not as leaves) in the tree before incorporating full nodes into the tree. The goal is to keep the number of sessions small for which a node forwards packets for.

Algorithm CLUSTER(s_i):

- (1) Let P be the sequence of partial nodes in \mathcal{R}_i that are sorted in the decreasing order of available bandwidth capacity. Let U be the set of full nodes in \mathcal{R}_i , sorted by decreasing available bandwidth capacity.
- (2) Determine the minimum m for which a tree T can be constructed connecting all nodes in \mathcal{R}_i where m nodes in U are non-leaf nodes and a possibly empty subset of nodes $P' \subseteq P$ as non-leaf nodes in which
 - $D_{\max}(T) \leq \Delta$
 - nodes with higher available bandwidth capacity appear at a lesser depth (closer to the source).
 - any set of partial nodes P'' that satisfies $P' \subset P'' \subseteq P$ cannot be used as non-leaf nodes in a tree T' that satisfies $D_{\max}(T') \leq \Delta$.
- (3) If such an m exists, build the tree from m nodes drawn from U and the remaining nodes drawn from the set P' . Otherwise, return that no such tree can be built.

We note that determining the minimum m for a given set P' can be done in time $O(|\mathcal{N}| \log |\mathcal{N}|)$. This follows from a result in work by Malouch *et al* [12] that proves that the minimum depth tree is the one where nodes with greater bandwidth

availability are placed closer to the source of the tree. Hence, it is sufficient to build the tree by first sorting nodes in order of decreasing fanout, and then attaching the nodes to the tree in this order to the node of minimal depth that contains available edges.

Figure 2 shows a simple example where two sessions are sharing a network of 8 nodes. The number assigned to each node in the figure is the aggregate bandwidth available to support additional connections, and we assume that the rate of each session is 1 unit.

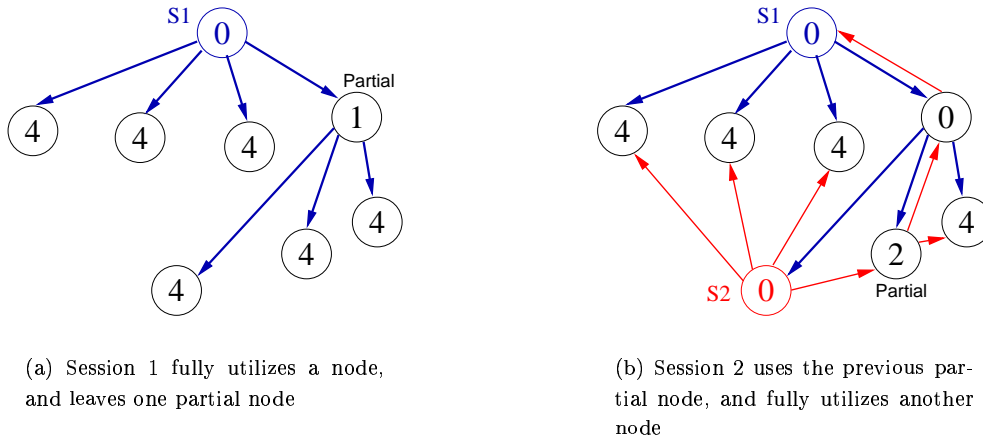


Figure 2: CLUSTER: an example with $|\mathcal{N}| = 8$, $\beta_n = 4 \forall n \in 1 \dots 8$, $\rho = 1$, $\Delta = 3$

3.2 Algorithm DISPERSE

Algorithm DISPERSE forms trees by utilizing as many nodes as possible for forwarding without exceeding the depth constraint. Barring a depth constraint, a chain of nodes is the preferred tree for DISPERSE (Figure 3).

Algorithm DISPERSE(s_i):

- (1) Let S be the sequence of nodes in \mathcal{R}_i sorted in order of decreasing available bandwidth capacity. Let c be the number of children that the first member of S can support for the session.
- (2) Let $d_{\min}(c)$ be the depth of a minimum-depth tree, $T(c)$ that can be built using no more than bandwidth c from any node in S . Again, we use results in [12]

to form this tree in time $O(|\mathcal{R}_i| \log |\mathcal{R}_i|)$. If $d_{\min}(c) > \Delta$, return that a tree cannot be constructed.

- (3) Form the tree $T(c-1)$ with depth $d_{\min}(c-1)$. If $d_{\min}(c-1) > \Delta$, return tree $T(c)$. Otherwise, decrement c and go to step (2).

Note that DISPERSE does more than merely “dispersing” the bandwidth in that it also tries to find the adequate fanout (c) that should be used to satisfy the depth constraint.

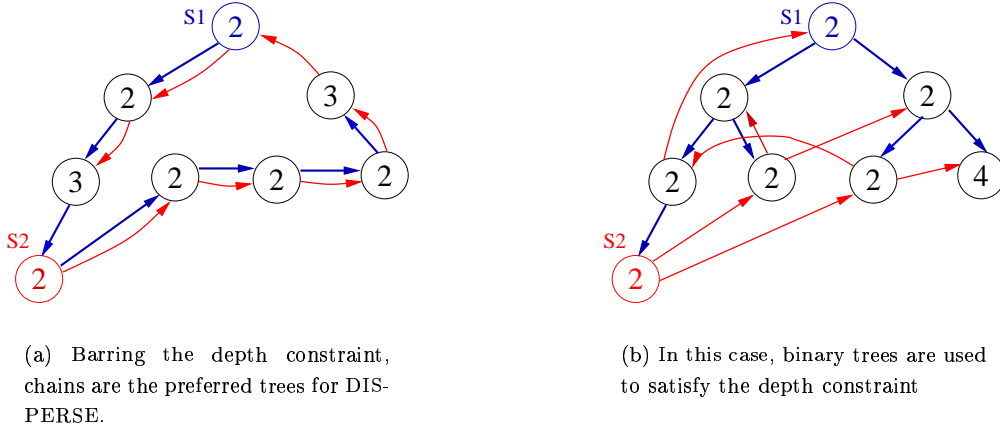


Figure 3: DISPERSE: an example with $|\mathcal{N}| = 8$, $\beta_n = 4 \forall n \in 1 \dots 8$, $\rho = 1$, $\Delta = 3$

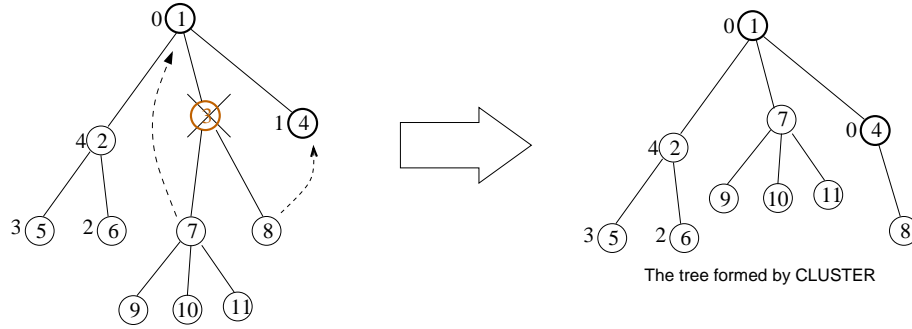
3.3 Algorithm Extensions

The algorithms described above have two shortcomings. The first shortcoming addresses the utilization of the node that will transmit data. This node may not have sufficient bandwidth to forward its own transmissions because its bandwidth is being utilized to forward transmissions of other sessions. We consider three possible actions that we refer to as *variants* that the network can take when faced with such a situation.

- no-swap: If s_i does not have enough outgoing capacity to initiate a session, the session is rejected.
- swap: s_i shifts some of its forwarding responsibility to another node that is participating in the same sessions as s_i that has the bandwidth capacity to

take on the additional bandwidth load. If no other node can be located that has sufficient spare capacity, the session to originate from s_i is dropped.

- reservation: A fixed amount of outgoing capacity is reserved at every node specifically for the transmission of a session that is initiated at that node. At other times, the bandwidth must remain unused.



The subtree of greatest depth is attached at the highest node: node 7 is attached to the root of the tree at height 0. Then node 8 is attached to a node at height 1. CLUSTER selects node 4 with the minimum available bandwidth.

Figure 4: The leave procedure

The second shortcoming addresses the *dynamics* of membership within sessions. In particular, it is likely that there will be applications in which session members arrive late and leave early. Late arrivers must be added to the multicast tree within the depth bounds. By leaving a session or all sessions, a node that forwards transmissions for the sessions will disconnect the tree for each session it leaves. This disconnection must be repaired.

Our attachment and disconnect procedures for the two algorithms attempt to mimic the traits exhibited by the original algorithm. For CLUSTER, a late joiner is attached to the node with minimum available bandwidth that is sufficient to support an additional transmission such that the joiner is added at or above the depth bound. If no such attachment point can be provided, then the late joiner is attached to a node beyond the depth bound that has sufficient bandwidth to support such a node, when one exists. The node is otherwise rejected. The reason why we proceed to attach nodes in situations where the depth bound is violated will be explained below.

When a node leaves the session, its children become the roots of detached subtrees, where the subtrees can have variable depth. These subtrees are then attached at the highest points at which there is available capacity to perform the attachment, with the subtrees of greatest depth being attached at the highest point. When there are several nodes at the same height to which an attachment can be made, the node with minimum available bandwidth is selected (Figure 4). A subtree is discarded when there is insufficient available bandwidth to connect the root of the subtree.

For DISPERSE, a late joiner is attached to the node with maximum available bandwidth that lies within the depth bound. If no such node exists, then a node with maximum available bandwidth that lies outside the depth bound is used. If no node has sufficient available bandwidth to forward transmissions to the late joiner, the late joiner is rejected. The algorithm to re-attach detached subtrees after a node leaves the session is the same as that for CLUSTER, except that the node with maximum available bandwidth is selected when there are several nodes at the same height to which an attachment can be made (in Figure 4, node 8 is attached to node 2 instead of node 4).

The tree formed as a result of nodes joining and leaving the session may be able to admit fewer nodes into the session beyond the depth bound than one constructed in a static setting. For this reason, the tree is periodically reconfigured by applying the CLUSTER algorithm or the DISPERSE algorithm to the set of nodes that are currently participating in the session. Hence, nodes that are attached beyond the depth bound can be brought within the depth bound during this reconfiguration. Those nodes that lie outside the depth bound after a reconfiguration are dropped from the session.

4 Homogeneous Network

In this section, we present a formal proof that algorithm CLUSTER with swap is optimal within the homogeneous setting. By homogeneous, we mean that $\mathcal{R}_s = \mathcal{N}$ for all $s \in \mathcal{S}$, where $\rho_s = \rho_{s'}$ for all $s, s' \in \mathcal{S}$, and where $\beta_n = \beta_{n'}$ for all pairs $n, n' \in \mathcal{N}$. We set $F = \lfloor \beta_n / \rho_s \rfloor$, the the number of session copies that each node can forward, and let Δ_{\min} be the minimum depth over all trees that can be constructed within the available bandwidth constraints that includes all receivers. We borrow the following notation from [12]: $D_n(T)$: is the depth of node n in tree T where $D_s(T) = 0$ for the source of the tree s and $D_n(T) \equiv D_{\Pi_n(T)}(T) + 1$, where $\Pi_n(T)$ is

the parent in the tree T of node n . $D_{\max}(T, S) = \max_{n \in S} D_n(T)$, i.e., the maximum depth in T of any nodes in $S \subseteq T$. For conciseness, we define $D_{\max}(T) \equiv D_{\max}(T, T)$.

Our theoretical results rely on results proven previously in [12]. We restate the two definitions and a Lemma presented in that work that will be useful here.

Definition 1 (Fanout Constraints) *A fanout constraint function (FCF), $f()$, is a function that maps each node $n \in \mathcal{N}$ to a non-negative integer. $f(n) = i$ implies that node n has sufficient bandwidth capabilities to forward session data to at most i other nodes.*

Definition 2 *Let $f()$ and $g()$ be two different FCFs in G . We say that $f()$ **majorizes** $g()$ and write $f() \succ g()$ if there exist two orderings of the nodes in G , $\alpha_1, \dots, \alpha_m$ and β_1, \dots, β_m such that all the following hold:*

- $f(\alpha_i) \geq f(\alpha_j)$ and $g(\beta_i) \geq g(\beta_j)$ for $i < j$.
- For all $j > 0$, $\sum_{i=1}^j f(\alpha_i) \geq \sum_{i=1}^j g(\beta_i)$.
- $\sum_{i=1}^m f(\alpha_i) = \sum_{i=1}^m g(\beta_i)$.

We say that a tree T_0 is *legally connected* with respect to FCF $f()$ if each node n in T_0 has no more than $f(n)$ children. We will make use of the majorization result, proved in [12]:

Lemma 1 *Let $f()$ and $g()$ be two different maximum fanout assignments in which $f() \succ g()$. If T_0 is a legally connected tree with respect to $g()$, then there exists a tree T_1 with the same set of nodes that is legally connected with respect to $f()$ where $D_{\max}(T_1) \leq D_{\max}(T_0)$.*

In other words, if the total number of edges extending from all nodes must remain fixed, the minimum depth of a tree is reduced by having some nodes with a very large number of edges and others with a very small number of edges than when all nodes have roughly the same number of edges.

We let T_k equal the k th tree that is constructed with source node s_k by Algorithm CLUSTER. We define the FCF, $f_k(n)$, to equal the number of edges that connect from node n to its children in tree T_k . During the running of the algorithm, we define the k th *configuration* of sessions, C_k to be the set containing the first k trees

$\{T_1, \dots, T_k\}$ to be constructed, with C_0 being the configuration containing no trees. We define $f_{\Omega(C_k)}(n)$ to be the number of *residual* or as-of-yet unused edges of a node n within configuration C_k (i.e., after the first k trees have been constructed). Because our setting is homogeneous, we have that $f_{\Omega(C_0)}(n) = F$ for all $n \in \mathcal{N}$. We define a partial node $p \in P$ to be *usable* under configuration C_k whenever $f_{\Omega(C_k)}(p) < F$ and it is possible to build a tree T' where $D_{\max}(T') \leq \Delta$ using node p and along with other nodes n that satisfy $f_{\Omega(C_k)}(n) = F$. In other words, it is feasible to use p to build tree T_{k+1} .

Lemma 2 *Under Algorithm CLUSTER, tree T_1 contains at most one non-leaf node n_0 for which $f_1(n_0) < F$. Letting m equal the number of non-leaf nodes in tree T_1 that satisfy $f_1(n) = F$, then remaining trees, T_k contain either m or $m + 1$ non-leaf nodes, where at most two nodes, n_1, n_2 exist for which $f_k(n_j) < F, j = 1, 2$. Last, under any configuration C_k , there is at most one usable node in P .*

Proof: The fact that T_1 contains at most one node n_0 for which $f_1(n_0) < F$ follows trivially from algorithm CLUSTER, given that all nodes initially have $f_{\Omega(C_0)}(n) = F$, and CLUSTER turns at most one node into a partially utilized node for each tree that it builds. We show the last property in the lemma by induction. Assume the result holds true after building k trees, $k \geq 1$. The algorithm first uses any usable nodes in P when constructing the $k + 1$ st tree. Furthermore, the algorithm does not introduce an additional non-leaf node $n \in U$ until all nodes already present in the tree have all edges utilized. Hence, in the end, at most one non-leaf node in the tree is partially utilized, completing the inductive step.

Last, we note that when no usable nodes are available in P in configuration C_k , tree T_{k+1} will be isomorphic to T_1 . If a usable node p exists in P within configuration C_k , the tree must still contain at least $m - 1$ nodes n for which $f_{k+1}(n) = F$ (or else T_1 could have been constructed using fewer than m nodes that satisfy this property). If these $m - 1$ nodes and p do not produce a sufficient number of edges to connect the remaining (leaf) nodes, then an m th node is brought in. If this is not sufficient, then an $m + 1$ st node is brought in. This will clearly be sufficient since m of these $m + 1$ nodes satisfy $f_{\Omega(C_k)}(n) = F$ and hence could be used to build a tree isomorphic to T_1 . In this case, node p and the last node to which children are assigned in are the only two nodes that have children in T_{k+1} that do not apply all F edges within this tree. ■

Lemma 3 *When $\Delta > \Delta_{\min}$, then Algorithm CLUSTER(s) builds the maximum number of trees possible with depth no greater than Δ .*

Proof: We begin by noting that since all senders transmit at the same rate, if a node's remaining available bandwidth is insufficient to be a non-leaf node in tree T_i , it will also be insufficient for all trees T_j for $j > i$. We prove inductively on i that Algorithm CLUSTER builds T_i of depth at most $\Delta_{\min} + 1$, that only one node remains in P after the building of each tree, and that this one node is usable. The inductive step holds for T_1 : straightforward application of Lemma 1 and 2 that gives us that T_1 builds a tree of depth Δ_{\min} .

To continue our proof by induction, we first note that a node $p \in P$ within configuration C_k is always usable in the building of tree T_{k+1} . This can be seen by considering the worst case, where p has only one available edge to support downstream nodes, i.e., $f_{\Omega(C_k)}(p) = 1$. Here, it is possible to build T_k within the depth constraint by having the source send to p , and then extend the rest of the tree from p . This subtree can be built isomorphic to T_1 with one leaf node removed (since one node p has already been added). Since T_1 has depth Δ_{\min} , the T_k described above has depth $\Delta_{\min} + 1$. For the case where p has more than one usable edge, these additional edges could be used to further decrease the depth of T_k .

The Lemma is then proved by considering the FCF of remaining edges, $f_{\Omega(C_k)}()$ after the algorithm has generated k trees, and let $g()$ represent the fanout function of remaining edges of other nodes from any other arbitrary construction of k trees. Since each tree utilizes $n - 1$ edges, the total number of edges remaining is the same in the two cases, i.e., $\sum_{n \in \mathcal{N}} f_{\Omega(C_k)}(n) = \sum_{n \in \mathcal{N}} g(n)$. We have $f_{\Omega(C_k)}() \succ g()$ for all $k \leq \lfloor \sum_{n \in \mathcal{N}} \beta_n / \rho_s(|\mathcal{N}| - 1) \rfloor$ from the fact that $f_{\Omega(C_k)}()$ contains at most one node in P . The proof is then completed using Lemma 1. ■

Similarly, we can prove the optimality of CLUSTER for the tightest depth bound Δ_{\min} . Also, using the above Lemma together with the interchange argument ¹, we could extend the result to a dynamic setting where session initiation and termination times occur in some arbitrary order.

¹see for example Liu et al. [11] for a comprehensive treatment on this technique

5 Queuing Model: A Stochastic Knapsack

We derive performance bounds on our algorithms by considering equivalent queuing systems. We map our system into a stochastic knapsack framework [16], and compute lower bounds on blocking probabilities based on the equivalent system.

Let \mathcal{M} be the set of all possible bandwidth requirements of Sessions in the system (i.e. product of transmission rate and session size-1). Let M be $|\mathcal{M}|$. The stochastic knapsack consists of C resource units to which objects from M classes arrive. Objects from class m are distinguished by their arrival rate, λ_m , mean holding time, $1/\mu_m$ and their size b_m (Figure 5).

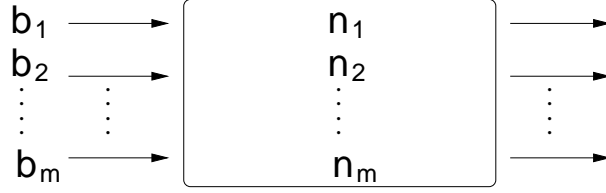


Figure 5: The stochastic knapsack of capacity $\sum_{n \in \mathcal{N}} \beta_n$

Let n_m denote the number of class- m objects in the knapsack. Then the total amount of resource utilized by the objects in the knapsack is $\mathbf{b} \cdot \mathbf{n}$, where $\mathbf{b} := (b_1, b_2, \dots, b_M)$ and $\mathbf{n} := (n_1, n_2, \dots, n_M)$. We define the process in terms of the state space of the different class- m objects, i.e. let

$$\mathcal{K} := \{\mathbf{n} \in \mathcal{I}^m : \mathbf{b} \cdot \mathbf{n} \leq C\}$$

The knapsack always admits an arriving object when there is sufficient room. More specifically, it admits an arriving class- m object if $b_m \leq C - \mathbf{b} \cdot \mathbf{n}$. Let \mathcal{K}_m be the subset of such states, i.e.

$$\mathcal{K}_m := \{\mathbf{n} \in \mathcal{K} : \mathbf{b} \cdot \mathbf{n} \leq C - b_m\}$$

. The blocking probability B_m for a class- m call under Poisson arrival assumption is then given by [16]

$$B_m = 1 - \frac{\sum_{\mathbf{n} \in \mathcal{K}_m} \prod_{j=1}^M \rho_j^{n_j} / n_j!}{\sum_{\mathbf{n} \in \mathcal{K}} \prod_{j=1}^M \rho_j^{n_j} / n_j!} \quad (1)$$

As a parsimonious metric to compare performance, we look at the weighted blocking probability w_b for each Stochastic Knapsack, where we weigh the blocking probability for a class- m call with the resource request of the call b_m , i.e. if $\mathbf{B} := (B_1, B_2, \dots, B_M)$, then

$$w_b = \frac{\mathbf{B} \cdot \mathbf{b}}{\sum_{m=1}^M b_m}$$

To compute the upper bound on the number of admitted sessions in our system (which correspondingly yields a lower bound on the blocking probability), we ignore the delay constraints and map the capacity C of the knapsack into $\sum_{n \in \mathcal{N}} \beta_n$, i.e. the aggregated bandwidth in the set of nodes. Thus, whenever there is available capacity in the knapsack (available bandwidth in the nodes), we'll assume the session is admitted irrespective of the delay constraints. The size of the object b_m is then simply the product of the size of the session (minus one) and the sending rate of the session. We also assume that the bandwidth requirements b_m are integers.

6 Simulation Results with Static Group Membership

In this section, we compare the performance of the algorithms proposed in Section 3 via an evaluation of their blocking probability, i.e., the fraction of sessions that must be turned away because the participating receivers do not have sufficient bandwidth capabilities to support the arriving session. To perform this evaluation, we use simulation. The number of nodes, $|\mathcal{N}|$ equals 100 in all simulation runs. Sessions arrive at random points at time, where the arrival times are described by a Poisson process with rate $\lambda = 1$. These sessions last for a time that is exponentially distributed with rate ρ , which we vary over the simulation runs. For each arriving session we randomly select one node to be the source. Here, we restrict each node to being the source of at most one session at any instant in time.

We now describe the various experiments individually and discuss the conclusions that we draw from each set of experiments. In all figures, we vary ρ , the expected lifetime of a session, along the x -axis. The y -axis indicates the blocking probability, and the various curves plot results for the various versions of the CLUSTER and DISPERSE algorithms we consider. The weights are set proportional to the rate and the size of the session, such that the penalty for dropping a session is proportional to the amount of bandwidth required by the session. For the case where all sessions are offered at the same rate and the same size, this weighted blocking probability simply reduces to the straightforward blocking probability.

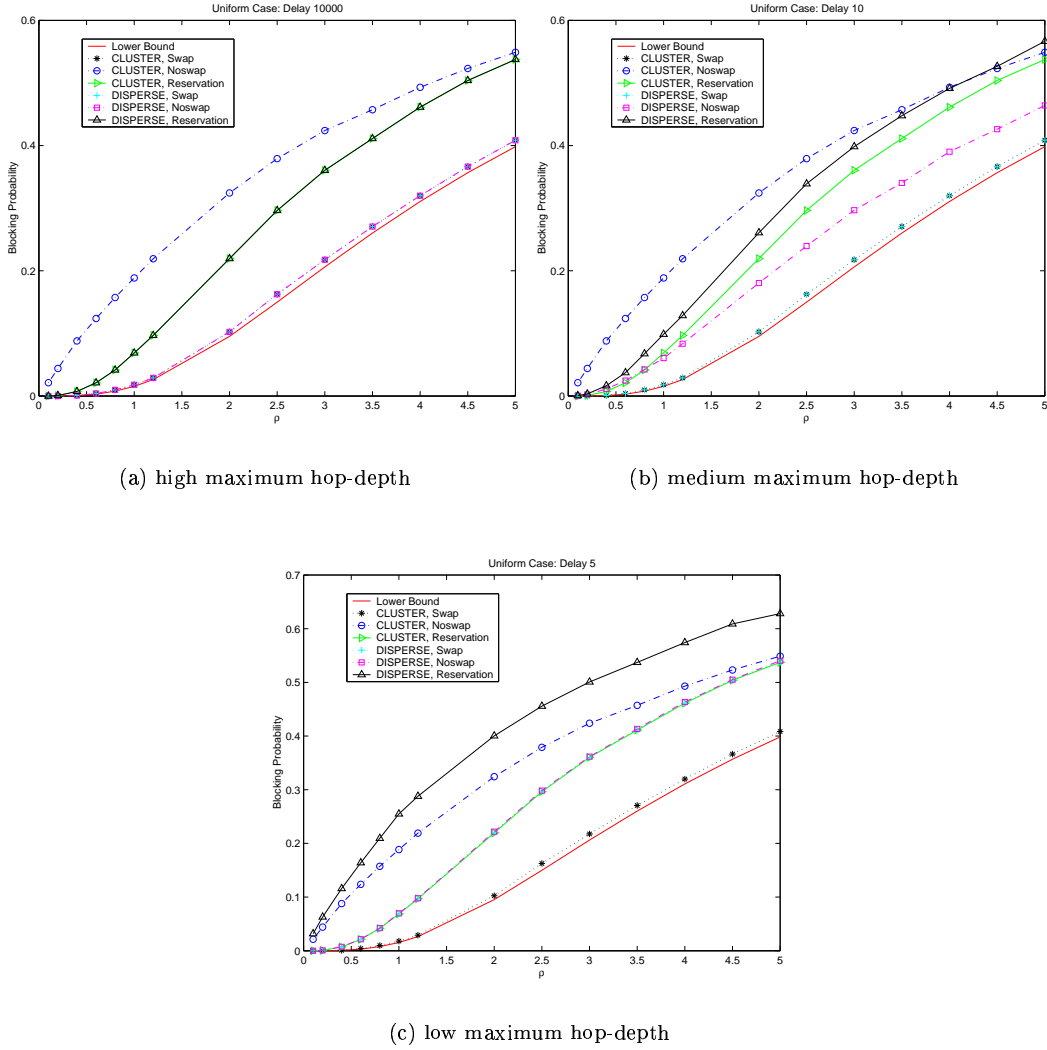


Figure 6: Homogeneous Setting

For each experiment, we plot results from three sets with different maximum hop-depth constraints: (a) high, (b) medium, and (c) low. We also plot the theoretical lower bound on the weighted blocking probability computed from the Stochastic Knapsack framework.

6.1 Experiment 1: Homogeneous case (Figure 6)

In this set of experiments, all receivers participate in all sessions. All session rates are fixed at 1 unit, and node capacities are fixed at 4 units. Since each session requires a set of 99 transmissions to reach all participants, a session consumes an aggregate bandwidth of 99 units.

From these experiments(Figure 6), we see that when there is a loose maximum hop-depth constraint, CLUSTER and DISPERSE with swap produce blocking probabilities that are close to the optimal, and that DISPERSE without swap has a significantly lower blocking probability than CLUSTER without swap. With a very tight maximum hop-depth constraint of 5, we observe that the blocking probability of CLUSTER with swap is close to the optimal, whereas the blocking probability of DISPERSE is unaffected by implementing swapping. CLUSTER without swap has a higher blocking probability than these other three. For both the algorithms and for all maximum hop-depths, reservation yields the highest blocking probability. A succinct description of the conclusions from these experiments is that for a homogeneous session configuration, CLUSTER should be used if swapping is permitted. Otherwise, DISPERSE yields a lower blocking probability.

6.2 Experiment 2: Variable Node Capacities (Figure 7)

In this set of simulations, the bandwidth capacity units of each node is selected uniformly at random from the set $\{3, 4, 5\}$. Session size is fixed at 100 and each session's rate is fixed at 1. Our conclusions from an examination of the plots in Figure 7 are for the most part similar to those reached in the case of a homogeneous network setting. Reservation again yields a higher blocking probability than the other variants. If swapping is enabled, CLUSTER and DISPERSE yield similar blocking probabilities. The one exception is the case where the maximum hop-depth is low. Here, using CLUSTER results in marginal improvements in blocking probability in comparison to DISPERSE. Again, the conclusion to be drawn is that if swapping is not permitted, DISPERSE yields a lower blocking probability. Otherwise, the choice of CLUSTER or DISPERSE is arbitrary since the blocking probabilities of the two

algorithms are virtually identical. The only difference is that the simulated blocking probabilities are significantly higher than the computed theoretical lower bound on the blocking probability.

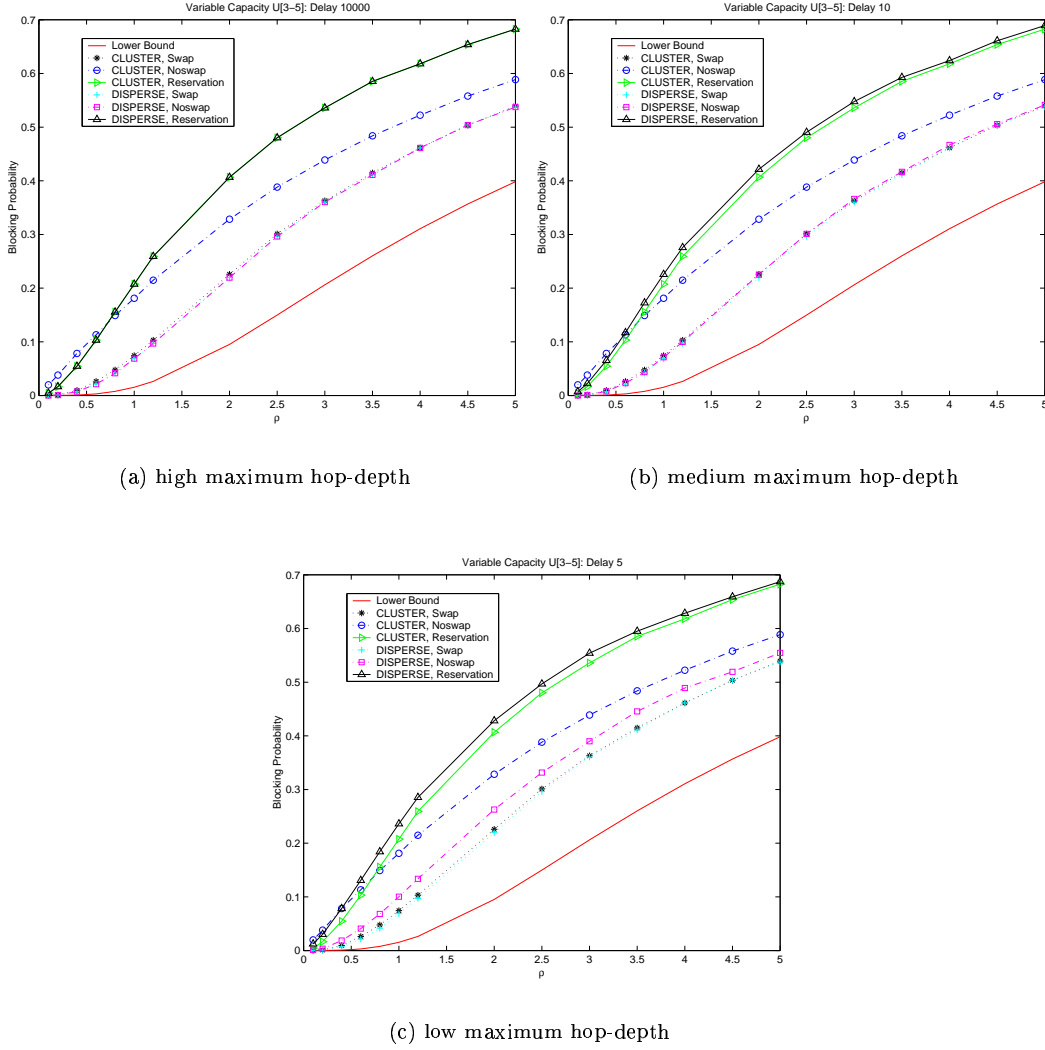


Figure 7: Heterogeneous setting, variable node capacity

6.3 Experiment 3: Variable Session Rates (Figure 8)

In this set of simulations, the rates of arriving sessions are selected uniformly at random from the set of integers ranging from 1 to 3. Node capacity is fixed at 4 and session size is fixed at 100. Examination of the plots in Figure 8 reveals that The weighted blocking probability that results from using DISPERSE meets the theoretical lower bound for large (unbounded) maximum hop-depths. For tighter maximum hop-depths, CLUSTER with swapping produces a slightly lower weighted blocking probability than DISPERSE, but the weighted blocking probability without swapping is significantly higher. Here, we conclude that when session depths can be unbounded, or when swapping is not permitted, DISPERSE should be used. Otherwise, CLUSTER yields moderate improvements in weighted blocking probability.

6.4 Experiment 4: Variable Session Sizes (Figure 9)

In this last set of simulations, the number of receivers that participate in the arriving session is selected uniformly at random from the set of integers ranging from 25 to 75. Node capacity is fixed at 4 units and session rates are fixed at 1 unit. Examination of the plots in Figure 9 reveals that with swapping, CLUSTER and DISPERSE yield approximately the same blocking probabilities, or CLUSTER's blocking probability is marginally lower. When swapping is not permitted, the blocking probability of CLUSTER is significantly higher than DISPERSE. We again conclude that if swapping is enabled, CLUSTER and DISPERSE yield almost identical blocking probabilities. When swapping is not enabled, DISPERSE should be used.

6.5 Reservation Policy

Unlike the telecommunications area in which reservation schemes lead to reductions in blocking of sessions, the reservation scheme analyzed here do not exhibit similar performance gains. This is in part due to the fact that the use of low-bandwidth end-users as forwarding agents required reservations of up to 25% of the user's bandwidth (i.e., one forwarding capacity unit). It will be interesting to evaluate these reservation schemes in environments in which the end users' bandwidth capabilities are larger such that the reservation threshold can be reduced to a smaller fraction on the order of 10%. In that case, we suspect that a reservation scheme can lead to significant decreases in blocking probabilities.

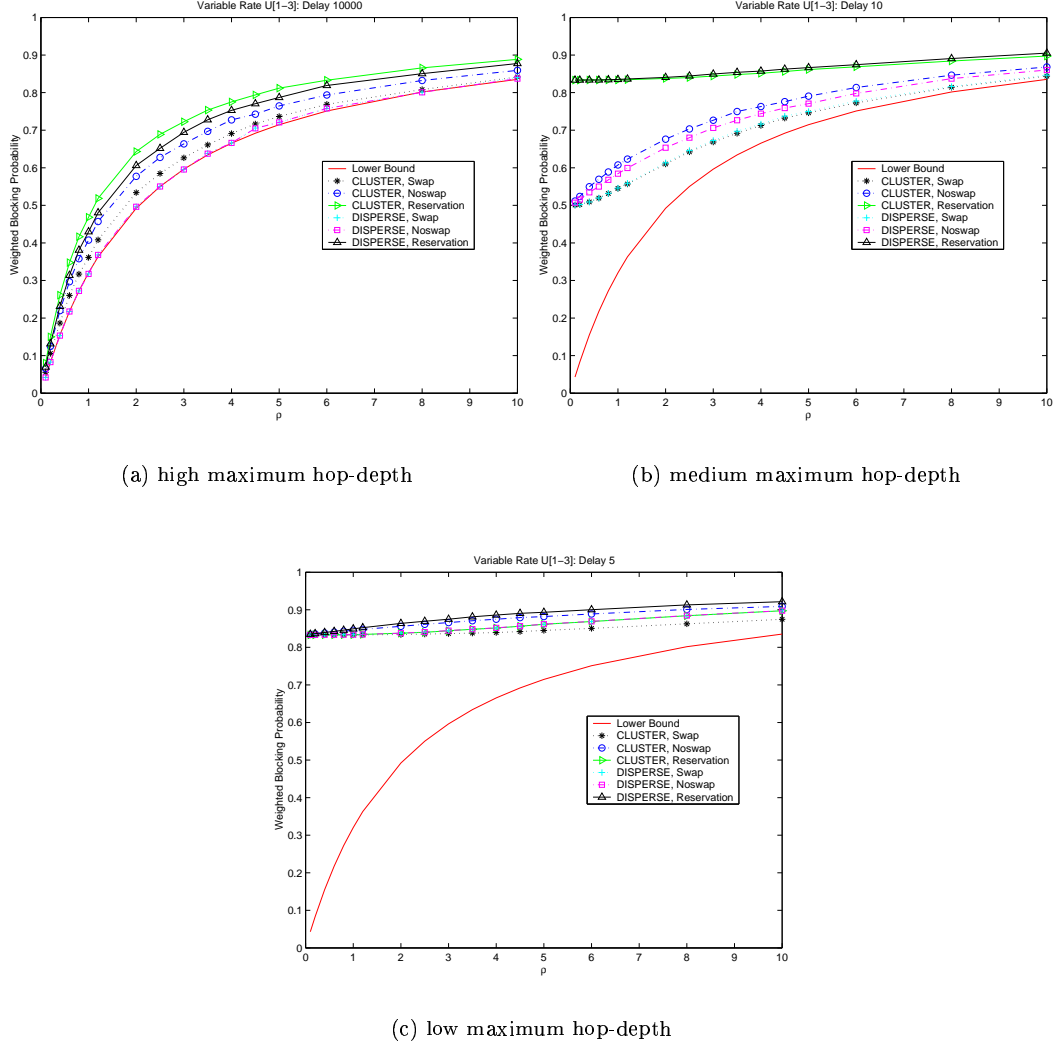


Figure 8: Heterogeneous setting, variable session rate

7 Simulations with Dynamic Group Membership

In this section we present results of simulations conducted with dynamic node membership. Individual nodes join and leave the sessions, and we use the dynamic ver-

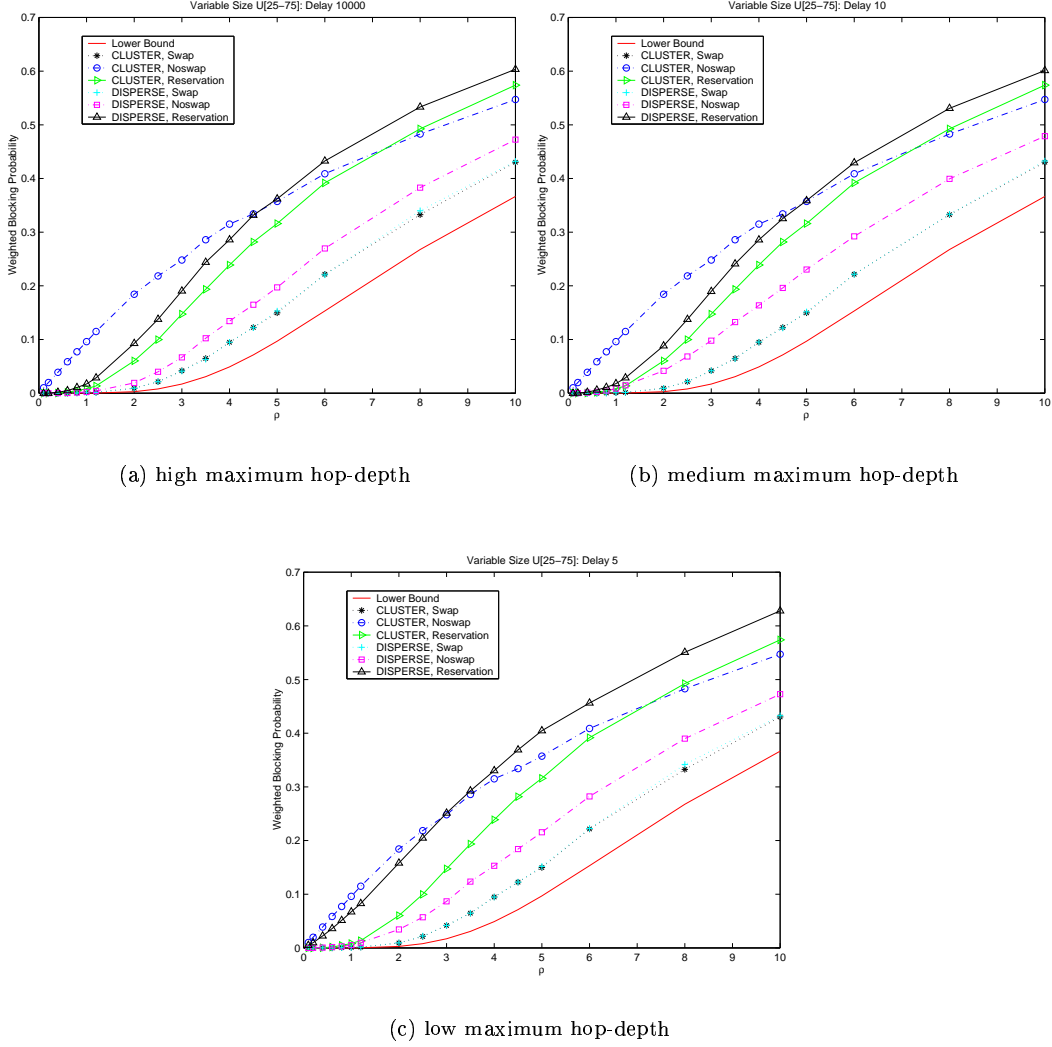


Figure 9: Heterogeneous setting, variable session size

sions of CLUSTER and DISPERSE described in the end of Section 3. We simulate two scenarios: in one, the number of *sessions* is fixed, but individual nodes join to and leave from this static set of sessions. In the other, the set of active sessions themselves as well as the membership within those sessions vary over time. For all

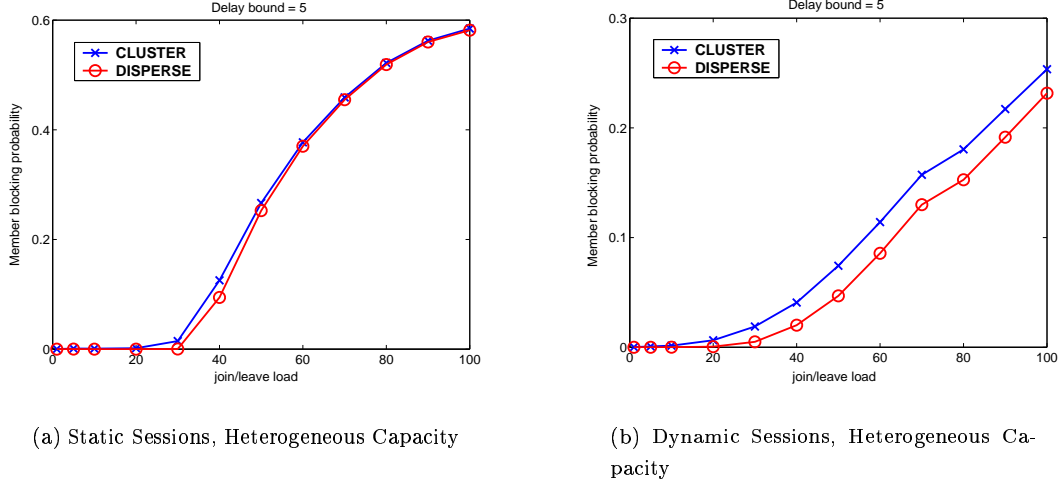


Figure 10: Dynamic join/leave experiments, Heterogeneous capacities

experiments, the maximum session size is 100 nodes, and each session's transmission utilizes one unit of bandwidth. For the static session case, the total number of sessions is fixed at 10, the join process is Poisson with rate 100 and the inter-update interval equal to 10. For the dynamic session case, the session arrival process is Poisson with rate 1, the session duration is exponentially distributed with a mean of 5 and the join rate within each session is again 100. The inter-update interval for the dynamic session case is much smaller, 0.05, as the topology in this scenario evolves much faster than in the static session case. For both scenarios, the member sojourn time is exponentially distributed, and we vary the mean along the x -axis, where the value indicated along the x -axis is the join rate divided by the leave rate. We plot the member blocking probability along the y -axis. We evaluate both the static and dynamic session case using homogeneous node capacities of 4, as well as heterogeneous node capacities uniformly distributed among $\{3, 4, 5\}$. Due to space constraints, we only present results where node capacities are heterogeneous and where there is a (tight) depth constraint of 5. We observe similar trends in the graphs for higher delay constraints and for homogeneous node capacities.

In Figure 10(a), we see that the blocking rates that result from applying CLUSTER and DISPERSE are similar, especially when the join rate divided by the leave rate is high. Intuitively, this is because under high loads, available bandwidth is

scarce, so that the typical tree formed is often the same, whether formed via CLUSTER or DISPERSE. For dynamic sessions, the blocking rates that result from applying DISPERSE are lower than those from applying CLUSTER (Figure 10(b)). Intuitively, we see two reasons why this is so. First, since in CLUSTER a node is either a leaf node, or (with high probability) completely utilized, a leave operation leads to a bursty process of re-attachment of the subtrees. Either there is no subtree to be re-attached, or there are a high number of nodes that must be re-attached, the latter requiring more bandwidth units. In contrast, with DISPERSE the available bandwidth is spread out among nodes. Hence, after a leave operation subtrees are more easily accommodated. Besides, the number of subtrees to re-attach is smaller than the one produced by CLUSTER even if the number of nodes to re-attach is the same. Indeed, in the best case, DISPERSE has to re-attach only one chain which needs one bandwidth unit.

Second, DISPERSE exploits better the fact that the depth is unbounded between reconfigurations. Since in the leave operation the subtrees are re-attached in order to balance the whole tree, a node that has joined the session beyond the depth bound, could be brought within the bound either when other nodes leave or when the reconfiguration is performed.

8 Conclusion

In this paper, we have explored the problem of building depth-bounded multicast trees for multiple sessions in networking systems where tree depth and a node's outgoing bandwidth constrain the permitted topology of the tree. We consider two algorithms that build overlay trees within these constraints, one tries to cluster a node's available bandwidth within a single tree, the other tries to disperse the available bandwidth among multiple trees. We derive a lower bound for the blocking probability of algorithms in this networking environment and compare the performance of our algorithms to this lower bound through simulation.

An interesting finding of our paper is that the clustering algorithm provably minimizes the session blocking rate in homogeneous network environments. On the other hand, we find through simulation that the dispersing algorithm exhibits a lower blocking rate in heterogeneous networking environments where session sizes, session rates or node capacities differ. Furthermore, the dispersing algorithm performs better in both homogeneous and heterogeneous settings in environments where session participants join and leave in the middle of a session.

In summary, our study indicates that it is more efficient to spread each node's forwarding capacity across sessions given the heterogeneous nature of real networking environments. Our on-going work consists in specifying the distributed algorithms for dynamically building the trees and determining how and when it is appropriate to perform the swapping operation.

References

- [1] Akamai Corporation. Internet Bottlenecks: The Case for Edge Delivery Services, 2000. Akamai whitepaper.
- [2] F. Bauer and A. Varma. Degree-constrained Multicasting in Point-to-point Networks. In *Proceedings of IEEE INFOCOM'95*, Boston, MA, March 1995.
- [3] Fred Bauer and Anujan Varma. Distributed Degree-Constrained Multicasting in Point-to-Point Networks. Technical Report UCSC-CRL-95-09, UCSC, 1995.
- [4] S. Bhattacharyya, D. Towsley, and J. Kurose. The Loss Path Multiplicity Problem for Multicast Congestion Control. In *Proceedings of IEEE INFOCOM'99*, New York, NY, March 1999.
- [5] Y. Chawathe, S. McCanne, and E. Brewer. RMX: Reliable Multicast for Heterogeneous Networks. In *Proceedings of IEEE INFOCOM'00*, Tel-Aviv, Israel, March 2000.
- [6] Y. Chu, S. Rao, S. Seshan, and H. Zhang. Enabling Conferencing Applications on the Internet Using an Overlay Multicast Architecture. In *Proceedings of ACM SIGCOMM'01*, San Diego, CA, August 2001.
- [7] Y. Chu, S. Rao, and H. Zhang. A Case for End System Multicast. In *Proceedings of ACM SIGMETRICS'00*, Santa Clara, CA, May 2000.
- [8] L. Gautier, C. Diot, and J. Kurose. End-to-end Transmission Control Mechanisms for Multiparty Interactive Applications in the Internet. In *Proceedings of IEEE INFOCOM'99*, New York, NY, March 1999.
- [9] Inktomi Corporation. The Inktomi Overlay Solution for Streaming Media Broadcasts. Inktomi whitepaper.

- [10] J. Jannotti, D. Gifford, K. Johnson, M. Kaashoek, and J. O'Toole. Overcast: Reliable Multicasting with an Overlay Network. In *Proceedings of USENIX*, San Diego, CA, October 2000.
- [11] Z. Liu, P. Nain, and D. Towsley. Sample Path Methods in the Control of Queues. *Queueing Systems, Special Issue on Optimal Control in Queueing Systems*, 21, 1995.
- [12] N. Malouch, Z. Liu, D. Rubenstein, and S. Sahu. A Graph Theoretic Approach to Bounding Delay in Proxy-Assisted End-system Multicast. In *Proceedings of the International Workshop on Quality of Service (IWQoS)*, Miami Beach, FL, May 2002.
- [13] S. McCanne, V. Jacobson, and M. Vetterli. Receiver Driven Layered Multicast. In *Proceedings of SIGCOMM'96*, Stanford, CA, August 1996.
- [14] D. Pendarakis, S. Shi, D. Verma, and M. Waldvogel. ALMI: An Application Level Multicast Infrastructure. In *3rd Usenix Symposium on Internet Technologies and systems (USITS)*, March 2001.
- [15] L. Rizzo. pgmcc: A TCP-friendly Single-Rate Multicast Congestion Control Scheme. In *Proceedings of ACM SIGCOMM'00*, Stockholm, Sweden, September 2000.
- [16] Keith W. Ross. *Multiservice Loss Models for Broadband Telecommunication Networks*. Springer-Verlag, 1995.
- [17] S. Shi and J. Turner. Dimensioning Server Access Bandwidth and Multicast Routing in Overlay Networks. In *Proceedings of NOSSDAV'01*, Port Jefferson, NY, June 2001.
- [18] S. Shi and J. Turner. Routing in Overlay Multicast Networks . In *Proceedings of IEEE INFOCOM'02*, New York, NY, June 2002.
- [19] B. Wang and J. Hou. Multicast Routing and its QoS Extension. *IEEE Network*, Jan/Feb 2000.
- [20] J. Widmer and M. Handley. Extending Equation-Based Congestion Control to Multicast Applications. In *Proceedings of ACM SIGCOMM'01*, San Diego, CA, August 2001.



Unité de recherche INRIA Sophia Antipolis
2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

Éditeur
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399